

# On Sharp Boundary Problem in Rule Based Expert Systems in the Medical Domain

*Oladipupo O. Olufunke, Covenant University, Nigeria*

*Uwadia O. Charles, University of Lagos, Nigeria*

*Ayo K. Charles, Covenant University, Nigeria*

---

## ABSTRACT

*Recently, the application of the conventional rule based expert system for disease risk determination in medical domains has increased. However, a major limitation to the effectiveness of the rule based expert system approach is the sharp boundary problem that leads to underestimation or overestimation of boundary cases, which ultimately affects the accuracy of their recommendation. In this paper, an expert driven approach is used to investigate the viability of a fuzzy expert system in the determination of risk associated with coronary heart disease with regards to the sharp boundary problem in rule based expert system.*

*Keywords: Coronary Heart Disease, Disease Risk Determination, Fuzzy Logic, Quantitative Binary Partition, Rule Based Expert System*

---

## 1. INTRODUCTION

The use of human expert knowledge in form of rules to solve real-world problems that normally would require human intelligence, known as rule based expert system, has played an important role in modern intelligent systems (Harleen & Siri, 2006). However, a major limitation to the effectiveness of this class of expert systems is the sharp boundary problem (SBP) which leads to underestimation or overestimation of boundary cases as a result of the quantitative attributes partitioning strategy; which conse-

quently affects the accuracy of the expert system (Verlinde, Cook, & Boute, 2006).

In the medical domain, the use of rule based expert system has increased greatly because of the scarcity of human experts in the domain and the availability of fast growing databases which could be used to model inferences and discover patterns in form of rules. In real live application, medical databases contain different kinds of attributes such as binary and quantitative attributes (Delgado, Marin, Sanchez, & Vila, 2001). Binary takes values from 0 or 1; for instance, a patients smoking status could be 'yes' or 'no'. Quantitative attributes that are categorical, numerical, or non-fractional in nature, take values from an ordered numerical

DOI: 10.4018/jhisi.2010070102

scale, often a subset of the real number (Kuok, Fu, & Wong, 1998). Quantitative attributes are very common in medical databases. For example heart disease patients can take age values between 20-79 years, result from laboratory test for systolic blood pressure level could take values within  $<120$  to  $\geq 160$ mm/Hg, cholesterol measures could be within the range of  $<160$  to  $\geq 280$ mg/dL.

In building an expert system, quantitative attributes need to be partitioned into ranges because of the very wide range of values defining their domain. There are several approaches to partitioning quantitative attributes as discussed in literature (Han & Kamber, 2001). The partitioning process is referred to as binning, that is an interval is considered as a "bin". The common binning strategies are: 1) Equiwidth binning, where the interval size of each bin is the same; 2) Equidepth binning, where each bin has approximately the same number of tuples assigned to it; and 3) Homogeneity-based binning, where bin size is determined so that the tuples in each bin are uniformly distributed. Also, there is the Distance based partitioning strategy, which seems most intuitive since it groups quantitative values that are closed together within the same interval (Han & Kamber, 2001). All of these partitioning strategies are subject to sharp boundary problem because of the classical set theory (Kuok, Fu, & Wong, 1998). However, to prevent this problem, in (Navruz & Serhat, 2007) fuzzy logic concept was introduced into a rule based expert system to determine coronary heart disease risk. The design gives the user the risk ratio and most of the experimented test data risk ratio from the fuzzy approach was reported to give relatively the same percentage risk as Adult Treatment Panel III (ATP III) calculation, which reflect the extent to which fuzzy concept was able prevent sharp boundary problem. In our work a comparative study was undergone to investigate the effect of SBP on quantitative binary partition strategy and fuzzy partition strategy in building a rule base expert system.

The starting point for fuzzy set theory (Zadeh, 1965) is that it is against intuition to

model vague concepts such as young and high by crisp intervals. For why would a person be considered as young while he is younger than 40, and on his 40th birthday suddenly loses this status? In the real sense, the transition between being youngAge and middleAge is not abrupt but gradual (Verlinde, Cook, & Boute, 2006). In fuzzy set theory, an element can belong to neighbouring sets each with set membership value in  $[0,1]$  depending on the type of membership function used. This value is assigned by each membership function associated with each fuzzy set. For attribute age and its domain  $D_{age}$ , the mapping of the membership function is  $\mu_{age}(x): D_{age} \rightarrow [0,1]$ . Fuzzy set is said to provide a smooth change between the boundaries (Kuok, Fu, & Wong, 1998). This is a very good argument for modeling vague concepts by fuzzy sets instead of crisp sets, as many researchers have already used this for the introduction of fuzzy logic to rule discovery (Cock, De Cornelis, & Kerre, 2003; Delgado, Marin, Sanchez, & Vila, 2003; Gyenesei, 2001). In Verlinde, Cook, and Boute (2006) an argument was actually raised against this in favour of binary partition strategy in association rule mining process. The argument was experimentally investigated using data driven approach. However, this argument could not be generalised since expert driven approach is yet to be considered (Verlinde, Cook, & Boute, 2006).

In this paper, Subject Matter Experts (SME), that is, medical doctors' knowledge (Ajith, 2005) is used as against the data driven approach used in (Verlinde, Cook, & Boute, 2006) to experimentally investigate the impact of SBP on rule based expert system to see if the introduction of fuzzy logic concept can have a significant impact on the accuracy of the rule based expert systems as introduced in (Navruz & Serhat, 2007). This investigation is concluded with a comparative analysis of quantitative binary partitioning expert system and fuzzy membership function partitions expert system with expert driven approach. This is applied to coronary heart disease (CHD) risk determination expert system based on Framingham risk

point scoring (Department of Health and Human Services, 2001; Bayliss, 2001).

The rest of this paper is organized as follows. In section 2, we review rule based expert system and Coronary heart disease as related to the Framingham risk scoring. Experimental results from our investigation are given in section 3, followed by conclusion in section 5.

## 2. LITERATURE REVIEW

### 2.1 Rule Based Expert System

Conventional rule-based expert systems use human expert knowledge to solve real-world problems that normally would require human intelligence (Ajith, 2005). Expert knowledge is often represented in the form of *rules* or as data within the computer. Depending upon the problem requirements, these rules and data can be recalled to solve problems. Rule based expert systems have played an important role in modern intelligent systems and their applications in strategic goal setting, planning, design, scheduling, fault monitoring, diagnosis and risk determination in medical domain and so on (Ajith, 2005). In generating the rule-base, several approaches have been introduced in literature (Department of Health and Human Services, 2001). The standard structure of a rule-base is such that given M dimensions where each dimension is partitioned into N subspaces, there exists up to  $N^M$  rules in the expert system (Aly & Vrana, 2006).

### 2.2 Fuzzy Logic

Fuzzy logic is the theory of fuzzy sets, sets that calibrate vagueness. It is a set of mathematical principles for knowledge representation based on degrees of membership. Unlike two-values boolean logic, fuzzy logic is a multi-valued logic, that allows intermediate values to be defined between conventional evaluations like true/false, yes/no, high/low, 1/0 etc. A fuzzy set is any set that allows its members to have different grades of membership in the interval  $[0,1]$ . For instance, let  $X$  be a collection of objects

and  $x$  a generic element of  $X$ , then a fuzzy set  $A$  in  $X$  is defined by function  $\mu_A(x)$  called the membership function of set  $A$ .

$$\mu_A(x): X \rightarrow \{0, 1\},$$

where  $\mu_A(x) = 1$  if  $x$  is totally in  $A$ ;

$\mu_A(x) = 0$  if  $x$  is not in  $A$ ;

$0 < \mu_A(x) < 1$  if  $x$  is partly in  $A$ .

To calibrate the membership value for every element of a fuzzy set, there are different membership functions, among which are Triangular MF (trimf), Trapezoidal MF (trapmf), Gaussian MF (gaussmf), Generalized Bell MF (gbellmf), etc. According to Zadeh, the nearer the value of  $\mu_A(x)$  to unity, the higher the grade of membership function of  $x$  in  $A$  (Zadeh, 1965). For example, using trapmf with the model in equation one (1) below, every element within  $10 \leq x \leq 60$  step 10 rage will be defined based on their membership values as:

$$\mu_{middleAge}(x) = \{10/0.0, 20/0.0, 30/0.5, 40/1, 50/1.0, 60/0.0\}$$

$$\mu_{middleAge}(x) = \begin{cases} \frac{(x-20)}{20} & 20 \leq x \leq 40 \\ 1 & 40 \leq x \leq 50 \\ \frac{(60-x)}{10} & 50 \leq x \leq 60 \end{cases} \quad (1)$$

### 2.3 CHD Risk Assessment

According to the Framingham CHD risks scoring the determinant factors for CHD risk are: age, total cholesterol, HDL cholesterol, systolic blood pressure, treatment for hypertension, and cigarette smoking. To determining 10-years risk the first step is to calculate the number of points for each risk factor. The total risk point sums the point for each determinant factor. The 10-years risk for myocardial infarction and coronary death is estimated from total points, and the

person is categorized according to absolute 10-years risk as indicated by Framingham risk assessment report (Department of Health and Human Services, 2001; Bayliss, 2001; Navruz & Serhat, 2007).

### 3. EXPERIMENTAL APPROACH

#### A. Data Sets

We conducted an experiment with 20 non-smoking men to determine their CHD risk according to the Framingham point score and to see the effect of SBP on their risk ratio. Using the standard structure of constructing a rule base such that given M dimensions where each dimension is partitioned into N subspaces, there exists up to  $N^M$  rules in an expert system rule base (Phayung, 2001). For this experiment we have 108 rules for the CHD risk determination expert system based on the number of input factors and the dimensions. Part of the developed rules is shown with Matlab rule interface in Figure 1. Table 1 shows the test data and the risk ratio for both binary (crisp) and fuzzy partitions. The graphical representation of the result is shown in Figure 4.

In order to analyse the impact of sharp boundary problem on CHD risk determination expert system, we need a comparison of crisp

and fuzzy partitioning of input and output quantitative attributes. For this experimental purpose, we assume that the determinant factors for CHD 10 years risk are age, cholesterol, high density lipoprotein (HDL) and systolic blood pressure. The four factors are quantitative attributes and they serve as the input for the rule construction. The output risk is also a quantitative attribute according to the medical expert used in the experiment and previous work from literature (Bayliss, 2001; Navruz & Serhat, 2007).

#### B. Quantitative Binary Expert System (QBES)

Quantitative Binary Expert System is developed using the binary partitioning strategy, whereby an element either belongs or not. For the input variables, we used equidepth partitioning method; this is because, for any specified number of intervals, equidepth partitioning minimizes the partial completeness level (Kuok, Fu, & Wong, 1998). And for the output variable, we used distance based partitioning strategy because it seems most intuitive, since it groups values that are close together within the same interval (Department of Health and Human Services, 2001) (Figure 2).

Figure 1. The Fuzzy rules

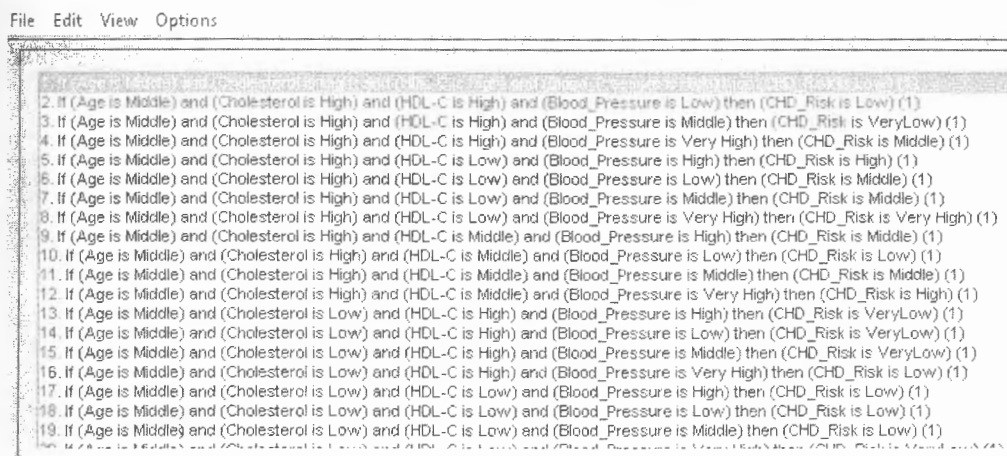


Table 1. ATP III, FES and QBES CHD % risk value according to 2+ risk factor CHD for non-smoking men

Patient no	Age	Cholesterol	HDL	Blood Pressure	ATP III	FES	QBES	ATP III CHD risk Linguistic value	FES CHD risk Linguistic value	QBES CHD risk Linguistic value
1	30	180	37	160	0	1.22	1.93	VeryLow	VeryLow	VeryLow
2	35	190	45	145	0	5.68	9.63	VeryLow	Low	Low
3	48	260	33	120	8	11.2	9.63	Low	Low	Low
4	57	300	67	110	8	9.68	9.63	Low	Low	Low
5	65	250	54	170	18	18.9	9.63	Middle	Middle	Low
6	75	290	25	135	30	31.8	32.5	VeryHigh	VeryHigh	VeryHigh
7	30	160	49	160	0	1.22	1.93	VeryLow	VeryLow	VeryLow
8	40	310	33	140	8	14.7	19.4	VeryLow	Middle	Middle
9	55	300	26	200	30	25.9	26.9	VeryHigh	High	High
10	60	230	39	110	11	11.3	9.63	Low	Low	Low
11	70	210	45	130	16	14.7	19.4	Middle	Middle	Middle
12	30	240	50	150	0	1.22	1.93	VeryLow	VeryLow	VeryLow
13	35	180	65	160	0	5.52	1.93	VeryLow	Low	VeryLow
14	45	300	47	155	9	14.7	19.4	Low	Middle	Middle
15	55	300	49	160	16	17.3	19.4	Middle	Middle	Middle
16	65	250	41	140	18	15.4	9.63	Middle	Middle	Low
17	70	260	38	190	30	26.7	32.5	VeryHigh	High	Middle
18	44	210	37	180	5	7	9.63	Low	Low	Low
19	55	150	30	200	11	17.3	19.4	Middle	Middle	Middle
20	66	150	26	200	28	24.5	19.4	High	High	Middle

For age, we have three partitions of young, middle and old. For cholesterol, we also have three partitions of Low, Normal and High. HDL is partitioned into three linguistic terms of Low, Middle and High. The Blood pressure is partitioned into four linguistic terms: Low, Middle, High, VeryHigh. Lastly for the output linguistic variable, CHD risk, we have 5 linguistic terms of VeryLow, Low, Middle, High, VeryHigh. The partition ranges and the Matlab representations are shown in Figure 3, Figure 4, Figure 5, Figure 6, and Figure 7.

In the experiment, we built quantitative binary expert system (QBES) based on binary partitioning strategy (either belong or not). Matlab fuzzy Tool box was used to simulate the expert system and the result is shown in Table 1. The Max-min operator of the Mandani fuzzy inference engine and centroid method for defuzzification process were used (Ajith, 2005). For example, for a non-smoking man of age 30, with Cholesterol 180 mg/dL, HDL-C 47 mg/dL, and bloodpressure 160mm/Hg, only rule 18 was fired and the calculated CHD risk is 9.63 as shown in Figure 8.

Figure 2. Input and Output variable partitioning (a) for Age, (b) Cholesterol, (c) HDL-C, (d) Blood pressure (e) CHD % risk

Age	Linguistics term	Cholesterol	Linguistics term	HDL	Linguistics term	Blood Pressure	Linguistics term	CHD Risk	Linguistics term
<20-39	Young	< 160-199	Low	< =39	Low	< 120	Low	<0-4	VeryLow
40-69	Middle	200-279	Normal	40-59	Middle	120-139	Middle	5-14	Low
>=70	Old	>=280	High	>=60	High	140-159	High	15-24	Middle
	(a)		(b)		(c)	>=160	VeryHigh	25-29	High
							(d)	>= 30	Very High
									(e)

Figure 3. Binary partition for age

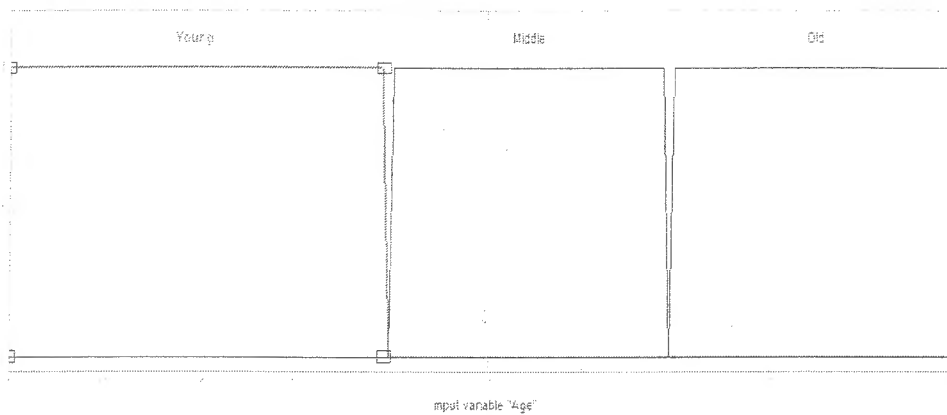
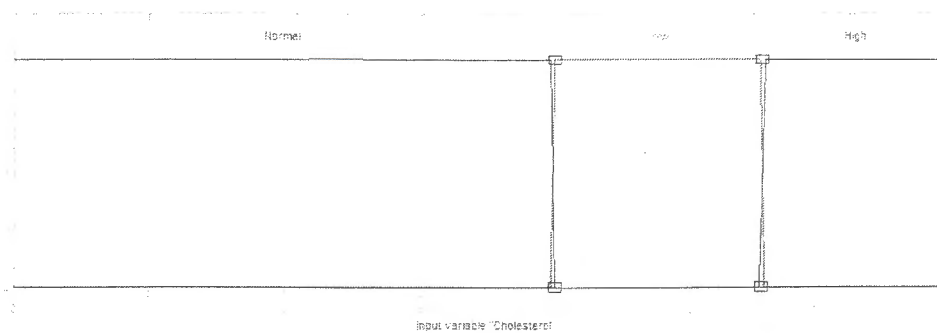


Figure 4. Binary partition for cholesterol



### C. Fuzzy Expert System (FES)

The fuzzy partitioning allows for overlapping of element within the neighboring linguistic terms which in turn prevents over estimation

of boundary values (Verlinde, Cook, & Boute, 2006). For input linguistic variables: age, cholesterol, HDC, blood pressure and output parameter CHD risk %, the fuzzy partitions are determined based on literature (Navruz &

Figure 5. Binary partition for HDL-C

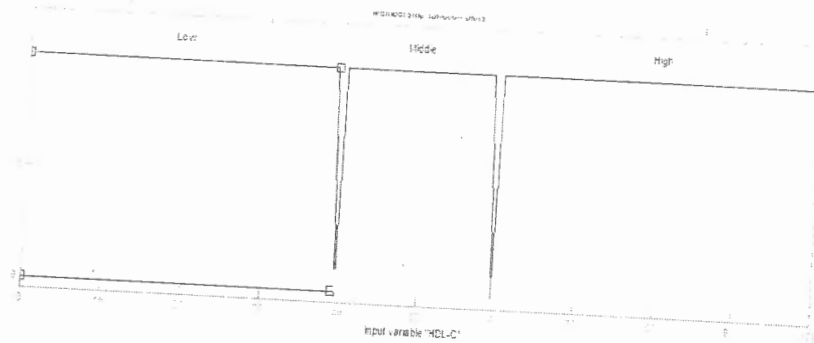


Figure 6. Binary partition for blood pressure

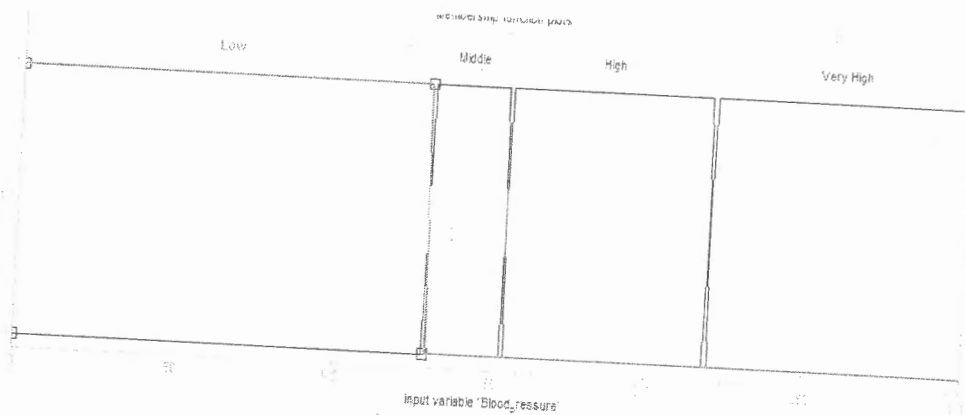


Figure 7. Binary partition for % CHD risk

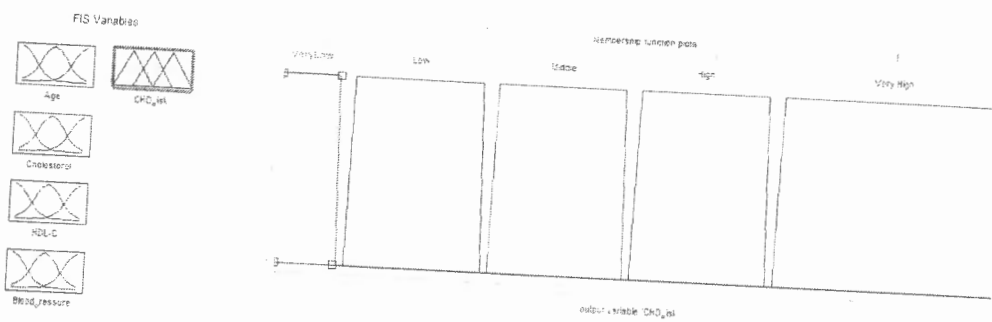
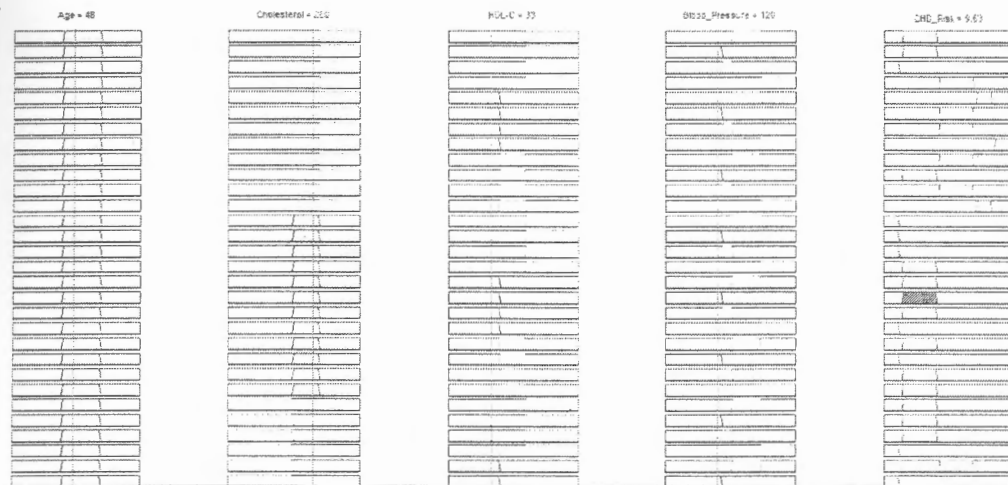


Figure 8. QBES CHD risk for the value Age=48, Cholesterol = 260, HDL-C=33, Bloodpressure = 120 with CHD risk % = 9.63.



Serhat, 2007) and the expert doctor. The trapezoidal membership function is used for fuzzy partitioning. For instance, age as an input is partitioned to youngAge, middleAge and old-Age as the linguistic values such that  $\forall x \in Age$ , the fuzzy membership models are:

$$\mu_{youngAge}(x) = \begin{cases} 1 & x \leq 20 \\ \frac{40-x}{20} & 20 \leq x \leq 40 \end{cases} \quad (2)$$

For other determinant factors the linguistic expressions are determined and their membership functions are represented with Matlab as shown in Figure 9, Figure 10, Figure 11, Figure 12 and Figure 13. The Max-min operator of the Mandani fuzzy inference engine and centroid method for defuzzification process was used. For example, a non-smoking man of age 30, with Cholesterol 180 mg/dL, HDL-C 47 mg/dL, blood pressure 160mm/Hg, rule 7,8,10,11,18,19,22 and 23 were fired and the calculated CHD risk is 11.2 as shown in Figure 13.

#### D. Comparing the Risk Ratio

From Table 1, it is observed that Fuzzy Expert System (FES) risk value varies as Adult Treatment Panel III (ATP III) risk based on the input variables, while Quantitative Binary Expert System (QBES) categorises different patients with different input values under the same risk. For instance, in considering patients 4, 5, 10, 16 and 18 from Table 1. The ATP III gives 8, 18, 11, 18 and 5 risk ratios, FES gives 9.68, 18.9, 11.3, 15.4, and 7 risk ratios respectively and QBES, gives 9.63 for all the patients. Categorically, this shows the effect of sharp boundary problem on the quantitative binary partitions; in that case, the five patients must have experienced either underestimation or overestimation of values as a result of the binary partitions.

From Table 1, Fuzzy Expert System (FES) risk ratios are considered 80% closer to Adult Treatment panel III(ATP III) as compared to Quantitative Binary Expert System (QBES). Also noticed, from Figure 7 and Figure 14, QBES fired only one rule to determine the output for the example cited because of the



Figure 9. The membership function for age

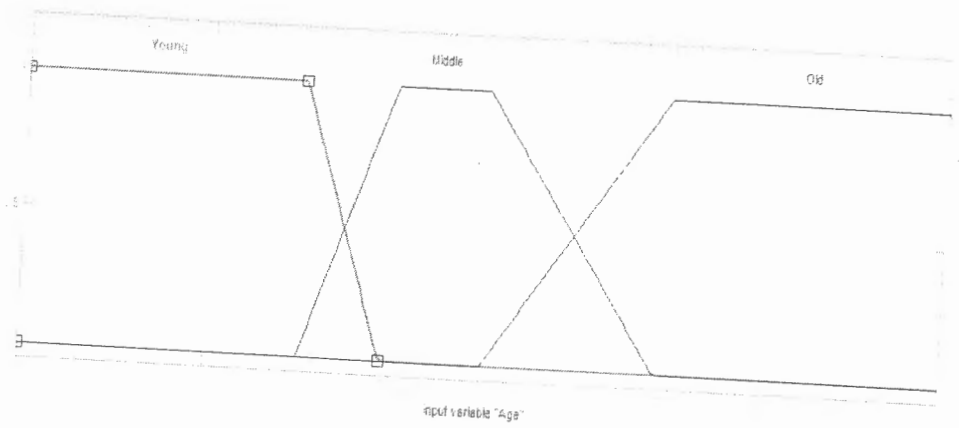


Figure 10. The membership function for cholesterol

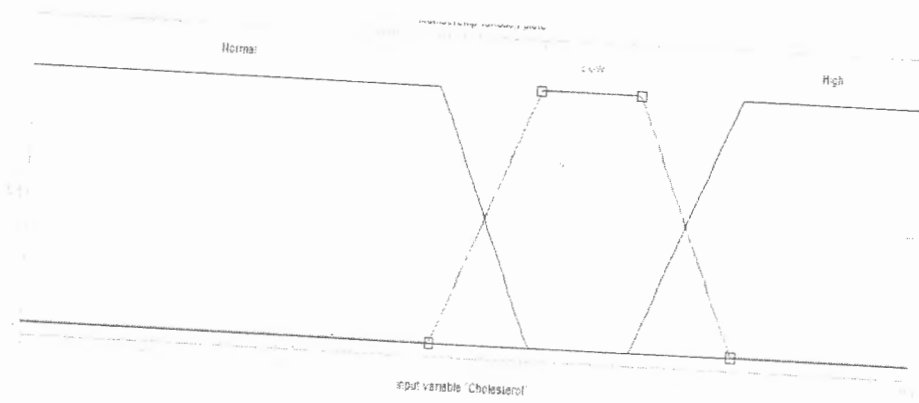


Figure 11. The membership function for (HDL-C)

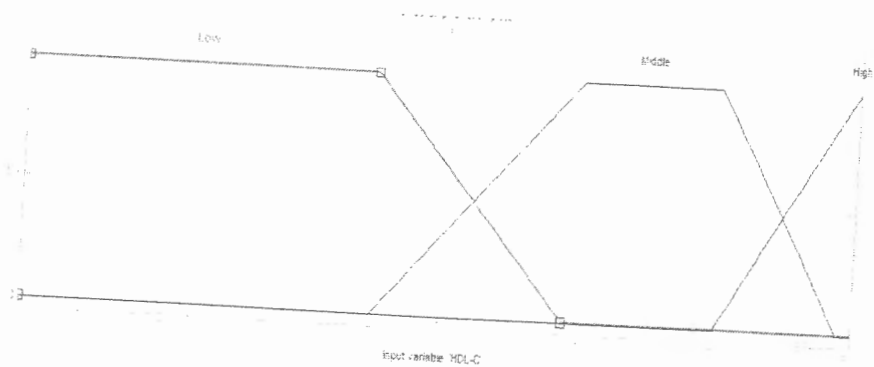


Figure 12. The membership function for blood pressure

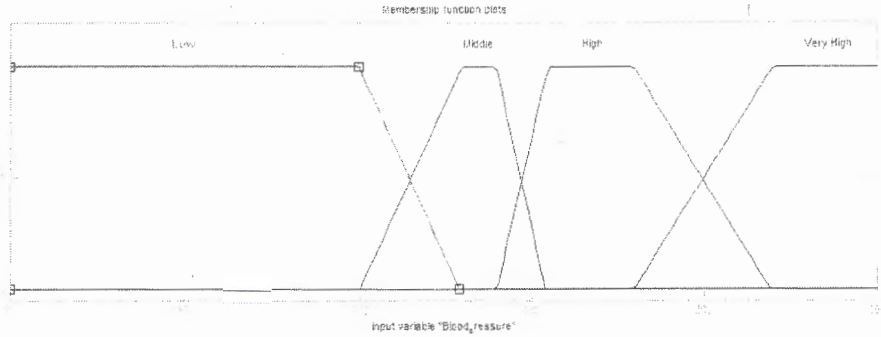


Figure 13. The membership function for CHD %risk

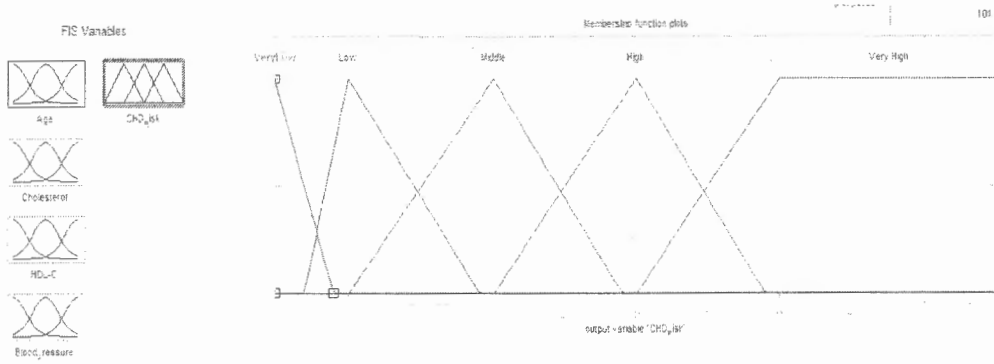
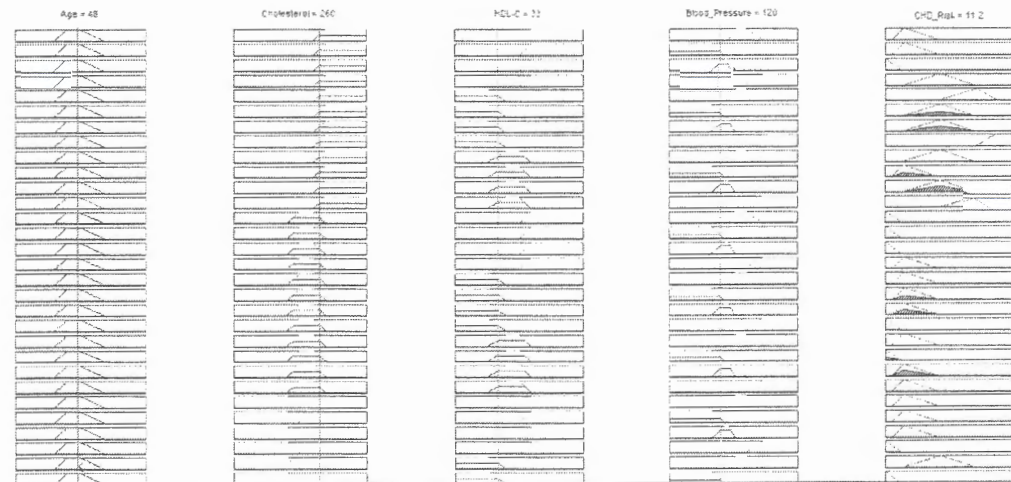


Figure 14. FES CHD risk for the value age=48, cholesterol = 260, HDL-C=33, blood pressure = 120 with CHD risk % = 11.2



binary partitioning process used. While FES, fired eight rule. This infers high percentage of closeness of FES and ATPIII outputs.

In order to get a better feel of the actual differences, Table 1 lists the 20 CHD patient records with ATP III 10 years risk according to Framingham report in Navruz and Serhat (2007) on 2+ risk factor for non-smoking men, FES risk values and QBES risk values. The chart for graphical overview is shown in Figure 15.

In Figure 16, linguistic values for CHD % risk: VeryLow, Low, Middle, High and Very-High are represented with 1, 2, 3, 4, and 5, respectively. The graph shows that in many instances the risk ratios fall under the same

linguistic value, but because of the domain under consideration, accuracy is paramount and essential. Therefore, the actual risk ratio is considered more important.

### 5. CONCLUSION

Fuzzy set theory is involved in expert systems for reasons which include suppression of unwanted problem that boundary element might cause (Verlinde, Cook, & Boute, 2006). This reason is convincingly proved when considering the risk outputs shown in Table 1. The differences between Fuzzy Expert System and

Figure 15. ATPIII, FES and QBES CHD % risk value diagramatic representation

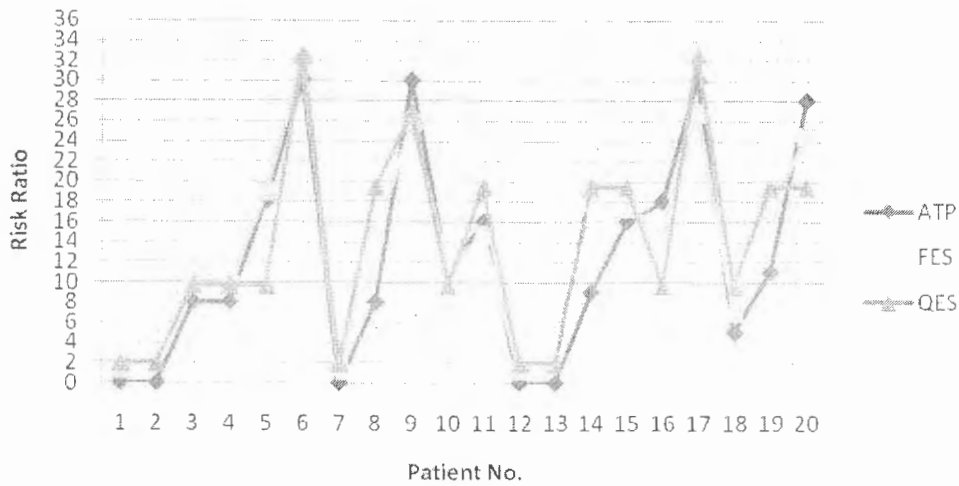
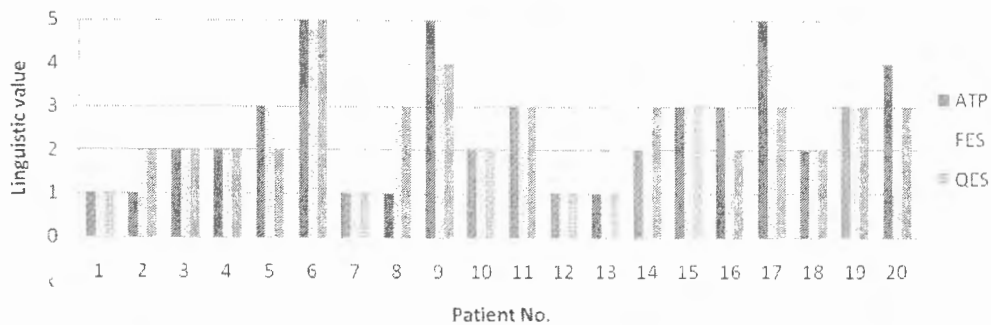


Figure 16. The linguistic % CHD risk diagramatic representation



Quantitative Binary Expert System risk ratios are significant because of the nature of medical domain, where lives are at stake. This is why a higher level of accuracy is required and the overestimation or underestimated of values due to the effect of Sharp Boundary Problem (SBP) cannot be tolerated.

Also, the introduction of expert knowledge in the partitioning process and the rule-base generation bring into a clear view the differences in the quantitative binary partitions and fuzzy membership partitions as against the opinion in Verlinde, Cook, and Boute (2006). These differences in turn affect the accuracy of a rule based expert system. Based on our investigation, we believe that the introduction of fuzzy logic in building rule based expert system minimizes the effect of SBP on boundary element.

In this experiment we have restricted ourselves to the use of few partitions for the determinant factors. In future work, it might be interesting to see the effect if the input variables fuzzy set (number of partition) is increased. For example if age could be further partitioned into VeryOld to make four partitions. Also, an investigation into enhancing the comprehensibility of FES through the use of interesting rules will be conducted.

## REFERENCES

- Ajith, A. (2005). *130: Rule-base Expert Systems. Handbook of measuring System Design*. New York: John Wiley & Sons.
- Aly, S., & Vrana, I. (2006). Toward efficient modeling of fuzzy expert systems: a survey. *AGRIC. ECON. - CZECH*, 52(10), 456-460.
- Barr, A., & Feigenbaum, E. A. (1982). The Handbook of Artificial Intelligence. *Information and Control*, 8(3), 338-358.
- Bayliss, J. (2001). Framingham risk score to predict 10 year absolute risk of CHD event west hertfordshire cardiology. In Wilson, P. W. F., (Eds.), *Prediction of coronary heart disease using risk factor categories*.
- Cock, M., De Cornelis, C., & Kerre, E. E. (2003). Fuzzy association rules: A two-sided approach. In *Proceedings of the Int. Conf. Fuzzy Information Processing- Theories and Applications*, Beijing, China (pp. 385-390).
- Delgado, M., Marin, N., Sánchez, D., & Vila, M.-A. (2001). Mining association rules with improved semantics in medical database. *Artificial Intelligence in Medicine*, 21, 241-245. doi:10.1016/S0933-3657(00)00092-0
- Delgado, M., Marin, N., Sánchez, D., & Vila, M.-A. (2003). Fuzzy association rules: General model and applications. *IEEE Transactions on Fuzzy Systems*, 11(2), 214-225. doi:10.1109/TFUZZ.2003.809896
- Department of Health and Human Services. (2001). *ATP III Guidelines At-A-Glance Quick Desk References (NIH publication No. 01-3305)*. Washington, DC: Department of Health and Human Services.
- Gyenesei, A. (2001). A fuzzy approach for mining quantitative association rules. *Acta Cybern.*, 15(2), 305-320.
- Han, J., & Kamber, M. (2001). *Data Mining Concepts and Techniques* (pp. 256-258). New York: Academic Press. ISBN1-55860-498-8
- Harleen, K., & Siri, K. W. (2006). Empirical Study on Applications of Data mining Techniques in Healthcare. *Journal of Computer Science*, 2(2), 194-200. ISSN 1549-3636
- Kuok, C. M., Fu, A. W.-C., & Wong, M. H. (1998). Mining fuzzy association rules in databases. *SIGMOD Record*, 27(1), 41-46. doi:10.1145/273244.273257
- Navruz, A., & Serhat, T. (2007). Design a Fuzzy Expert System for Determining of Coronary Heart Disease Risk. In *Proceedings of the International Conference on Computer Systems and Technologies- compSysTech'07*.
- Phayung, M. (2001). Quantitative measures of a Fuzzy Expert System. In *Proceedings of the IEEE Neural Network Council Student summer Research*.
- Srikant, R., & Agrawal, R. (1996). *Mining Quantitative Association Rules in Large Relational Tables*.
- Verlinde, H., Cook, M. E., & Boute, R. (2006). Fuzzy Versus Quantative Association Rules: A fair Data-Driven Comparision. *IEE Transactions on Systems, Man, and Cybernetics-Part B: Cybernetics*, 36(3).
- Zadeh, L. A. (1965). Fuzzy sets. *Information and Control*, 8(3), 338-358. doi:10.1016/S0019-9958(65)90241-X

*Oladipupo O. Olufunke holds a B.Sc., M.Sc. in Computer Science and she is currently a Ph.D. student in Covenant University. Her current research interests include data mining, fuzzy logic and expert system. She is a member of the Nigerian Computer Society (NCS), and Computer Professional Registration Council of Nigeria (CPN).*

*Uwadia O. Charles holds a B.Sc., M.Sc. and Ph.D. in Computer Science. His research interests include Software Engineering. He is the present President of the Nigerian Computer Society (NCS), and Computer Professional Registration Council of Nigeria (CPN). He is a Professor of Computer Science in the University of Lagos, Nigeria, Africa.*

*Ayo K. Charles holds a B.Sc., M.Sc. and Ph.D. in Computer Science. His research interests include mobile computing, Internet programming, e-business and government, and object oriented design and development. He is a member of the Nigerian Computer Society (NCS), and Computer Professional Registration Council of Nigeria (CPN). He is currently an Associate Professor of Computer Science in Covenant University, Ota, Ogun State, Nigeria, Africa.*